



Optimization Algorithms for Parameter Estimation and Data Reconciliation

L. T. Biegler
Chemical Engineering Department
Carnegie Mellon University
Pittsburgh, PA

Chemical
Engineering

Parameter Estimation and Data Reconciliation

- I. Motivation and Formulation
 - Problem Formulation
 - Derivation of Objectives
- II. Unconstrained Problems
 - Linear Least Squares
 - Gauss-Newton Methods
 - Trust Region Methods
- III. Constrained Least Squares
 - Problem Formulation
 - Specialization to SQP
 - EVM Problems
- IV. Statistical Inference
 - Linear Problems
 - Nonlinear Problems
- V. Introduction to Data Reconciliation
 - Optimization Problem
 - Gross Error Detection
- VI. Contaminated Normal
 - Derivation
 - Examples
- VII. Robust Statistics (M-Estimator)
 - Huber (Fair function)
 - Hamp1 (Redescending function)
 - Examples
- VIII. Mixed Integer Approaches
 - Formulations
 - Comparisons

2

Parameter Estimation References

- Bard, Y., *Nonlinear Parameter Estimation*, Academic Press, New York (1974)
- Boggs, P. T., R. H. Byrd, and R. B. Schnabel (1987), "A Stable and Efficient Algorithm For Nonlinear Orthogonal Distance Regression," *SIAM J. Sci. Stat. Comput.*, 8(6):1052-1078. (ODRPACK)
- Caracotsios, M., PhD Thesis, University of Wisconsin (1985) (GREG)
- Dennis, J.E., Gay, D.M., and Welsch, R.E. (1981), *ACM Trans. Math. Software*, Vol. 7, No. 3. (NL2SOL)
- Lawson, C.L., and Hanson, R.J. (1974), **Solving Least Squares Problems**, Prentice-Hall, Englewood Cliffs, N.J.
- More, J.J. (1978), *Lecture Notes In Mathematics No. 630*, G.A. Watson (ed.), Springer-Verlag, Berlin (MINPACK)
- Nocedal, J. and S. Wright, *Numerical Optimization*, Springer Verlag, Berlin (1998)
- Tjoa, I-B and L.T. Biegler, "Simultaneous Solution and Optimization Strategies for Parameter Estimation of Differential-Algebraic Equation Systems," *I&EC Research*, 30, p. 376 (1991).
- Tjoa, I-B and L. T. Biegler, "A Reduced Successive Quadratic Programming Strategy for Errors-in-Variables Estimation," *Computers and Chemical Engineering* 16, 6, p. 523 (1992)

3

Data Reconciliation References

Texts:

- Shankar Narasimhan, Cornelius Jordache, *Data Reconciliation and Gross Error Detection: An Intelligent Use of Process Data*, Gulf Pub Co, 2000
- J.A. Romagnoli, M. Sanchez, *Data processing and reconciliation for chemical process operation*, Academic Press International, ISBN 0-12-594460-8 (2000).

Papers:

- Tjoa, I-B and L.T. Biegler, "Simultaneous Strategies for Data Reconciliation and Gross Error Detection of Nonlinear Systems," *Computers and Chemical Engineering* 15, 10, p. 679 (1991)
- Albuquerque, J. S., and L. T. Biegler, "Gross Error Detection and Variable Classification in Dynamic Systems," *AIChE Journal*, 42, 10, p. 2841 (1996)
- Arora, N. and L. T. Biegler, "Redescending estimators for data reconciliation and parameter estimation," *Computers and Chemical Engineering*, 25, p. 1585 (2001)

and citations therein

4

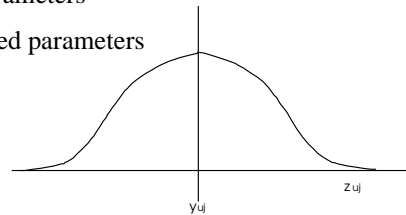
Maximum Likelihood Derivation - 1

Motivation:

- fit models to data to find 'optimal' parameters
- determine levels of confidence of fitted parameters
- evaluate suitability of models

Derivation of Objective Function

Let $z_{uj} = y_{uj}(\theta) + \epsilon_{uj}$
 z_{uj} - j th component of data for u th experiment
 y_{uj} - corresponding (correct) model value
 ϵ_{uj} - observation error following probability distribution function (pdf), $p(z)$
 θ - adjustable parameters in model



What is $p(z)$? For a scalar z , how is it derived?

5

Maximum Likelihood Derivation - 2

Assume the following moment information:

$$p(z) \geq 0, \int_{-\infty}^{\infty} p(z) dz = 1, \int_{-\infty}^{\infty} z p(z) dz = \eta, \int_{-\infty}^{\infty} (z - \eta)^2 p(z) dz = \sigma^2$$

Define measure of information (Shannon, 1948) and find the distribution function that maximizes the information assuming only the moment information:

$$\text{Max } I(p) = E(\log p) = \int_{-\infty}^{\infty} \log(p(z)) p(z) dz$$

$$\text{s.t. } p(z) \geq 0, \int_{-\infty}^{\infty} p(z) dz = 1$$

$$\int_{-\infty}^{\infty} z p(z) dz = \eta, \int_{-\infty}^{\infty} (z - \eta)^2 p(z) dz = \sigma^2$$

Problem can be solved analytically to yield: $p(z) = \frac{1}{(\sqrt{2\pi})\sigma} \exp(-(z - \eta)^2 / (2\sigma^2))$

If z is an m -vector for a single experiment, this can be extended to a joint multivariable distribution to give:

$$p(z) = (2\pi)^{-m/2} \det(V)^{-1/2} \exp(-1/2(z - \eta)^T (V)^{-1} (z - \eta))$$

where V is a covariance matrix defined by: $V = \int (z - \eta)(z - \eta)^T p(z) dz$

6

Maximum Likelihood Derivation - 3

Consider n multiple experiments (with index u). Each experiment has a mean η_u , covariance V_u and experimental error distribution, ϵ_u

$$p(\epsilon_u) = (2\pi)^{-m/2} \det(V)^{-1/2} \exp(-1/2(\epsilon_u)^T (V)^{-1}(\epsilon_u))$$

and the joint probability distribution is given by:

$$\prod_{u=1}^n p(\epsilon_u) = (2\pi)^{-mn/2} \left[\prod_{u=1}^n \det(V_u)^{-1/2} \right] \exp(-1/2 \sum_{u=1}^n (\epsilon_u)^T (V_u)^{-1}(\epsilon_u))$$

This distribution now needs to be converted into an objective function that 'maximizes information' about our data.

Let's make the following assumptions:

- Replace distributional errors, ϵ_u by the actual residuals, $e_u = (z_u - y_u(\theta))$
- Experiments u are independent and V_u is the same for all experiments, $V = E(\epsilon_u \epsilon_u^T)$
- Define likelihood function $L(\theta) = \prod p(e_u)$, and maximize this function (or its log).

This leads to the general form:

$$\log L(\theta) = -(nm/2) \log(2\pi) - (n/2) \log(\det(V)) - 1/2 \sum_{u=1}^n (e_u)^T (V)^{-1} (e_u)$$

7

Maximum Likelihood Objective Functions

Specialize objectives based on what we know about the error distributions.

Define moment matrix:

$$M(\theta) = \sum_{u=1}^n (e_u)(e_u)^T, \quad \text{Tr}(V^{-1}M(\theta)) = \sum_{u=1}^n (e_u)^T (V)^{-1} (e_u)$$

and $\log L(\theta) = -(nm/2) \log(2\pi) - (n/2) \log(\det(V)) - 1/2 \sum_{u=1}^n \text{Tr}(V^{-1}M(\theta))$

Since first two terms do not contain θ , we simply minimize $\text{Tr}(V^{-1}M(\theta))$

Special cases:

- **Ordinary Least Squares:** V is known, all component errors e_{uj} have same distribution and are independent of each other, i.e., $V = v I$

$$\text{Min Tr}(M(\theta)) = \sum_u \sum_j e_{uj}^2$$

- **Simple Weighted Least Squares:** V is known and diagonal, all component errors e_{uj} are independent of each other, i.e., $V = \text{diag}\{\sigma_j^2\}$

$$\text{Min Tr}(V^{-1}M(\theta)) = \sum_u \sum_j e_{uj}^2 / \sigma_j^2$$

- **Weighted Least Squares:** V is known but general, all component errors e_{uj} depend on each other:

$$\text{Min Tr}(V^{-1}M(\theta)) = \sum_u \sum_j e_u^T V^{-1} e_u$$

8

Maximum Likelihood – Unknown Covariance

Maximize $L(\theta)$ wrt V as well as θ .

$$\log L(\theta) = -(nm/2)\log(2\pi) - (n/2)\log(\det(V)) - 1/2 \sum_{u=1}^n \text{Tr}(V^{-1}M(\theta))$$

Assuming V is symmetric and nonsingular, we get:

$$\frac{\partial \log L(\theta)}{\partial V} = -(n/2)V^{-1} - 1/2 \sum_{u=1}^n (V^{-1}M(\theta)V^{-1}) = 0$$

which leads to: $V^* = 1/n M(\theta)$

Substitution into $L(\theta)$ leads to:

$$\log L(\theta) = -(nm/2)\log(2\pi) - (n/2)\log(n) - 1/2mn - n/2 \log \det(M(\theta))$$

so we minimize: $n/2 \log \det(M(\theta))$

Notes:

- If the structure of V is known, we can further specialize the objective for unknown covariance (e.g., diagonal covariance, same covariance)
- V^* is biased but can be corrected by using $V'' = n/(n-m) V^*$

9

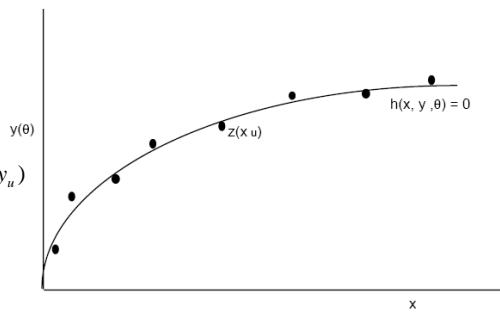
Unconstrained Least Squares

Basic Problem:

$$\text{Min } \Phi(\theta, y, z) = 1/2 \sum_{u=1}^n (z_u - y_u)^T W_u (z_u - y_u)$$

$$\text{s.t. } h_u(x_u, y_u, \theta) = 0$$

$$a \leq \theta \leq b$$



- W_u can be chosen to be $(V_u)^{-1}$
- $h_u(x_u, y_u, \theta) = 0$ is a model of the system
- x_u are fixed independent variables
- Bounds are not expected to be active

Consider the case where y_u is an explicit function of θ ,

$$h_u(x_u, y_u, \theta) = 0 \implies y_u = f_u(x_u, \theta)$$

$$\text{and } \text{Min } \Phi(\theta, y, z) = 1/2 \sum_{u=1}^n (z_u - f_u(\theta))^T W_u (z_u - f_u(\theta))$$

10

Least Squares – Solution Methods

Linear Least Squares

Model is given by: $y_u = A_u \theta + b_u$

$$\text{Min } \Phi(\theta, z) = 1/2 \sum_{u=1}^n (z_u - A_u \theta - b_u)^T W_u (z_u - A_u \theta - b_u)$$

From $\nabla_{\theta} \Phi(\theta, z) = 0$, we get the normal equation:

$$\sum_{u=1}^n (A_u^T W_u A_u) \theta = \sum_{u=1}^n A_u^T W_u (z_u - b_u)$$

and this leads to: $\theta = \left(\sum_{u=1}^n (A_u^T W_u A_u) \right)^{-1} \sum_{u=1}^n A_u^T W_u (z_u - b_u)$

Note: A better way to solve this linear system is to do a QR factorization on a concatenation (over u) of $(W_u)^{1/2} A_u$.

Least Squares – Solution Methods

Nonlinear Least Squares

Model is given by: $y_u = f_u(x_u, \theta)$

$$\text{Min } \Phi(\theta, z) = 1/2 \sum_{u=1}^n (z_u - f_u(\theta))^T W_u (z_u - f_u(\theta))$$

And we solve this unconstrained problem with Newton's method

From Taylor series expansion: $\nabla_{\theta} \Phi(\theta^k, z) + \nabla_{\theta\theta} \Phi(\theta^k, z) \Delta \theta = 0$

$$\nabla_{\theta} \Phi(\theta^k, z) = - \sum_{u=1}^n J_u^T W_u (z_u - f_u(\theta^k)), \quad J_u = \nabla_{\theta} f_u^T$$

$$\nabla_{\theta\theta} \Phi(\theta^k, z) = \sum_{u=1}^n (J_u^T W_u J_u + R_u), \quad \{R_u\}_{ij} = - \nabla_{\theta_i \theta_j} f_u^T W_u (z_u - f_u(\theta^k))$$

Now assume that $(z_u - f_u(\theta))$ is nearly zero and therefore R_u is nearly zero. Then the Hessian simplifies to: $\nabla_{\theta\theta} \Phi(\theta^k, z) = \sum_{u=1}^n J_u^T W_u J_u$

and we have the **Gauss-Newton Method**:

$$\Delta \theta = \left(\sum_{u=1}^n (J_u^T W_u J_u) \right)^{-1} \sum_{u=1}^n J_u^T W_u (z_u - f_u(\theta^k))$$

Note: A better way to solve this linear system is to do a QR factorization on a concatenation (over u) of $(W_u)^{1/2} J_u$.

Globalization of Gauss-Newton Method

To ensure convergence from poor starting points:

Line search method

Choose $\alpha \in (0, 1]$ so that a sufficient decrease is found for $\Phi(\theta)$ with: $\theta^{k+1} = \theta^k + \alpha \Delta\theta$.

This will converge to a stationary point ($\nabla_{\theta}\Phi = 0$) as long as

$\sum_{u=1}^n J_u^T W_u J_u$ is sufficiently positive definite. What if singular?

Add λI to Hessian to get the **Levenberg-Marquardt method**.

$$\Delta\theta = \left(\sum_{u=1}^n (J_u^T W_u J_u) + \lambda I \right)^{-1} \sum_{u=1}^n J_u^T W_u (z_u - f_u(\theta^k))$$

How should λ be adjusted?

13

Unconstrained Problems: TR Motivation

Take: $Min \Phi(\theta)$

Model Problem: $Min \Phi(\theta_k) + \nabla\Phi(\theta_k)^T \Delta\theta + \frac{1}{2} \Delta\theta^T \nabla^2\Phi(\theta_k) \Delta\theta$

$$\frac{1}{2} \|\Delta\theta\|^2 \leq \frac{1}{2} \Delta^2$$

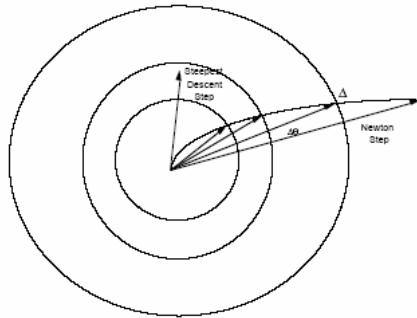
Optimality Conditions: Leads to Levenberg-Marquardt Method

$$\left\{ \begin{array}{l} \nabla\Phi(\theta_k) + \nabla^2\Phi(\theta_k)^T \Delta\theta + \lambda\Delta\theta = 0 \\ \|\Delta\theta\|^2 \leq \Delta^2, \lambda \geq 0 \\ \lambda(\|\Delta\theta\|^2 - \Delta^2)/2 = 0 \end{array} \right\} \implies \Delta\theta = -(\nabla^2\Phi(\theta_k) + \lambda I)^{-1} \nabla\Phi(\theta_k)$$

How do we choose λ or Δ ?

14

Extremes of Trust Region Method



For given Δ , solve for λ directly

$$1/\|\Delta\theta(\lambda)\| - 1/\Delta = 0$$

$\lambda = 0, \Delta$ large: $\Delta\theta$ is the Gauss-Newton Step

$\lambda \rightarrow \infty, \Delta = 0$: $\Delta\theta$ is a small step in the steepest descent direction

Trust Region Methods - guaranteed to converge to local optima, with much weaker assumptions than line search methods.

15

Trust region approach in MINPACK (More', 1980)

Choose λ so that $\|\Delta\theta\| \leq \Delta$, comparing $\tau = \text{ared}/\text{pred}$

- actual reduction (ared): $\Phi(\theta^k) - \Phi(\theta^k + \Delta\theta)$
- predicted reduction (pred): $\nabla_{\theta\theta}\Phi^T\Delta\theta + 1/2 \Delta\theta^T\nabla_{\theta\theta}\Phi \Delta\theta$

At iteration k:

- Calculate λ corresponding to Δ , and calculate $\Delta\theta$.
 - Evaluate (ared) and (pred).
 - Define $\tau = \text{ared}/\text{pred}$
 - a) If $\rho_1 < \tau \leq \rho_0$, $\Delta = m_1 \Delta$
 - b) If $\rho_2 \leq \tau < \rho_1$ or $\rho_0 < \tau$, $\Delta = \Delta$
 - b) If $\tau < \rho_2$, $\Delta = \Delta/m_1$
 - c) If $\tau < \rho_3$ reset $\Delta\theta = 0$
- Set, $\theta^{k+1} = \theta^k + \Delta\theta$

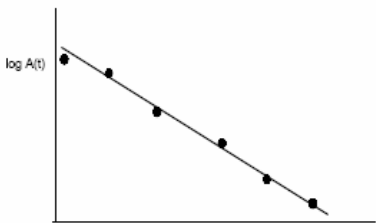
Typical values for the parameters:

$$m_1 = m_2 = 2 \text{ and } \rho_1 = 0.75, \rho_2 = 0.25, \rho_3 = 0, \rho_0 = 2.$$

16

Chemical Engineering

Why is trust region required for parameter estimation?



First order reaction: $A(t) = A(0) \exp(-k t)$; $k = k_0 \exp(-E/RT)$
 Data available only at one temperature.

Results:

- Nonunique parameter estimates, k_0 , E
- Singular Hessian
- Can be due to poor model and/or poor data
- Trust region methods will ensure convergence to “some” solution,
- Postoptimality analysis will establish nonuniqueness, insensitivity

17

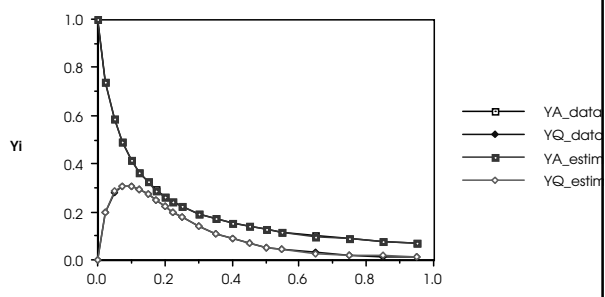
Chemical Engineering

Example: Catalytic Cracking of Gasoil (Tjoa, 1991)

k_1 k_2 k_3
 $A \rightarrow Q \rightarrow S, A \rightarrow S$

$$y_A' = -(k_1 + k_3) y_A^2$$

$$y_Q' = -k_1 y_A^2 - k_2 y_Q$$

$$y_A(0) = 1, y_Q(0) = 0$$


ODEs can be solved to yield explicit form: $y = f(\theta, t)$

Apply Trust Region method (GREG):

$(k_1, k_2, k_3)_0 = (6, 4, 1)$
 $(k_1, k_2, k_3)_* = (11.95, 7.99, 2.02)$
 $(k_1, k_2, k_3)_{\text{true}} = (12, 8, 2)$
 Converges in 5 iterations (11 function calls)

18

Convergence Rates for Gauss-Newton Methods

Small Residuals at Solution (Good Model Fit):

$$e_{u^*} = 0, R_{u^*} = 0,$$

$$\nabla_{\theta\theta} \Phi(\theta^*, z) = \sum_{u=1}^n J_u^T W_u J_u$$

- Gauss-Newton method is quadratically convergent.
- Trust region will be inactive if Hessian is nonsingular
- L-M is also quadratically convergent for unique θ^* .

Large Residuals at Solution (Poor Model Fit):

$$e_{u^*} \neq 0, R_{u^*} \neq 0,$$

$$\nabla_{\theta\theta} \Phi(\theta^*, z) = \sum_{u=1}^n J_u^T W_u J_u + R_u$$

- Gauss-Newton method is linearly convergent.
- Trust region may not be inactive if Hessian is nonsingular
- L-M is also linearly convergent for unique θ^* .

19

Hybrid Methods (Gauss and quasi-Newton)

Quasi-Newton Methods

- DFP and BFGS Methods apply secant formula, symmetry and positive definiteness of Hessian
- Do not exploit structure of least squares problem

Dennis, Gay and Welsh (1981) - NL2SOL

- Approximates true Hessian as G-N Hessian is known
- Specialized, self-scaling Q-N method developed that approximates R_u
- Incorporates Trust Region Approach of More'
- Leads to superlinear convergence

Fletcher and Xu (1987)

- Applies specialized Q-N method to approximate R_u
- Uses a switching rule to monitor if there are small or large residuals

$$\tau^k = (\Phi(\theta^k) - \Phi(\theta^{k+1})) / \Phi(\theta^k)$$

Large residuals: $\lim_{k \rightarrow \infty} \tau^k = 0$, Use specialized Q-N update if $\tau^k \leq \epsilon$

Small residuals: $\lim_{k \rightarrow \infty} \tau^k = 1$, Use specialized G-N formula if $\tau^k > \epsilon$

(Choose $\epsilon \sim 0.1$)

20

Constrained Least Squares

Motivation:

- Model cannot be reformulated as $y = f(\theta)$
- Too expensive to converge $h_u(y_u, \theta) = 0$ for each parameter value

Basic Formulation:

$$\begin{aligned} \text{Min } \Phi(\theta, y, z) &= 1/2 \sum_{u=1}^n (z_u - y_u)^T W_u (z_u - y_u) \\ \text{s.t. } h_u(y_u, \theta) &= 0 \\ a &\leq \theta \leq b \end{aligned}$$

- Any NLP method can be used to solve this problem
- SQP can be tailored to take advantage of special form of Φ allows for tailored algorithm.
- Leads to faster algorithm than standard SQP with BFGS updates

21

Optimization Strategy: SQP method

Let $x^T = [\theta^T, y^T]$ and consider QP subproblems for SQP:

$$\begin{aligned} \text{min } \nabla\Phi(x^k)^T d + 1/2 d^T B d \\ \text{s.t. } h(x^k) + \nabla h(x^k)^T d &= 0 \\ x^L \leq x^k + d \leq x^U \end{aligned}$$

First order necessary conditions

$$\begin{bmatrix} B & \nabla h \\ \nabla h^T & 0 \end{bmatrix} \begin{bmatrix} d \\ v \end{bmatrix} = - \begin{bmatrix} \nabla\Phi \\ h \end{bmatrix}$$

Problems:

- How to deal with a larger QP problems
- How to calculate the Hessian

Strategies:

- Use Range and Null space Decomposition strategy to decompose the search direction into:
 - Null space movement
 - Range space movement
- Use a hybrid Gauss-Newton and BFGS update formula
- Analogy to unconstrained approaches

22

Range and Null Space Decomposition

Define linear QP system as: $M s = f$, to give:

$$\begin{bmatrix} B & \nabla h \\ \nabla h^T & 0 \end{bmatrix} \begin{bmatrix} d \\ v \end{bmatrix} = - \begin{bmatrix} \nabla \Phi \\ h \end{bmatrix}$$

and select an $n \times n$ nonsingular matrix: $H = [Y \mid Z]$, where $\nabla h^T Z = 0$.

- Z, Y are null & range space bases for the linearized equalities
- Search direction with range (p_Y) and null space (p_Z) components:

$$d = Y p_Y + Z p_Z, \quad Y^T = [0 \mid I] \quad Z^T = [I \mid N^T C^{-T}]$$

Now defining $X = \text{diag} [[Y \mid Z], I]$, we can consider the equivalent system $X^T M X z = X^T f$ (with $X z = s$) as:

$$\begin{bmatrix} Y^T B Y & Y^T B Z & Y^T \nabla h \\ Z^T B Y & Z^T B Z & 0 \\ \nabla h^T Y & 0 & 0 \end{bmatrix} \begin{bmatrix} p_Y \\ p_Z \\ v \end{bmatrix} = - \begin{bmatrix} Y^T \nabla \Phi \\ Z^T \nabla \Phi \\ h \end{bmatrix}$$

Standard assumptions: set $Y^T B Y = 0$ and $Y^T B Z = 0$

23

Structure of Least Squares Hessian

$$B^k = \begin{bmatrix} \nabla_{\theta\theta} L & \nabla_{\theta y} L \\ \nabla_{y\theta} L & \nabla_{yy} L \end{bmatrix} \quad \begin{aligned} \nabla_{\theta\theta} L &= \sum_{j=1}^m v_j \nabla_{\theta\theta} h_j & \nabla_{\theta y} L &= \sum_{j=1}^m v_j \nabla_{\theta y} h_j \\ \nabla_{y\theta} L &= \sum_{j=1}^m v_j \nabla_{y\theta} h_j & \nabla_{yy} L &= \nabla_{yy} \Phi + \sum_{j=1}^m v_j \nabla_{yy} h_j \end{aligned}$$

where $L(\theta, y, v) = \Phi(y) + v^T h(\theta, y)$

KKT multipliers (based on first order estimates) are given by:

$$v = - (Y^T \nabla h)^{-1} Y^T \nabla \Phi = - (Y^T \nabla h)^{-1} Y^T [0 \mid \sum_u e_u^T W_u]^T$$

Assumption: If the residuals are *small*, then at convergence $e_u \approx 0 \Rightarrow v \approx 0$

The Hessian becomes
$$B^{G-N} = \begin{bmatrix} 0 & 0 \\ 0 & (W_u) \end{bmatrix}$$

\Rightarrow *Newton-like* rate of convergence

24

Least Squares Hybrid SQP Method

Motivation: Choose best Hessian approximation for different problem types

Strategies: Develop a switching rule to decide if Q-N or G-N approximation should be made for B^k . (Fletcher and Xu, 1987)

- Define merit function: $L^*(x^k) = \Phi(x^k) + v^T h(x^k) + 1/2 \gamma \|h(x^k)\|^2$
- Uses a switching rule to monitor if there are small or large residuals

$$\tau^k = (L^*(x^k) - L^*(x^{k+1})) / L^*(x^k)$$

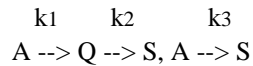
Large residuals: $\lim_{k \rightarrow \infty} \tau_k = 0$, Use specialized Q-N for $Z^T B Z$ if $\tau_k \leq \epsilon$

Small residuals: $\lim_{k \rightarrow \infty} \tau_k = 1$, Use $Z^T B^{GN} Z$, $Z^T B^{GN} Y_{p_Y}$ formula, $\tau_k > \epsilon$

Choose $\epsilon \sim 0.2$

25

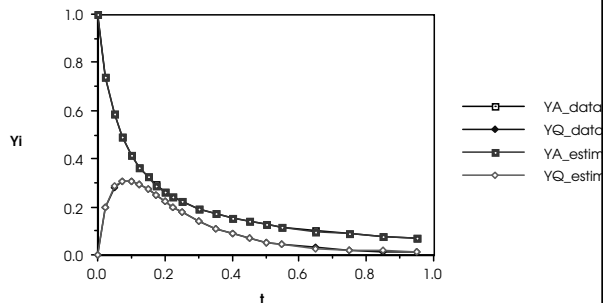
Example: Catalytic Cracking of Gasoil (Tjoa, 1991)



$$\begin{aligned}
 y_A' &= -(k_1 + k_3) y_A^2 \\
 y_Q' &= -k_1 y_A^2 - k_2 y_Q \\
 y_A(0) &= 1, y_Q(0) = 0
 \end{aligned}$$

number of ODEs: 2
 number of parameters: 3
 discretized ODEs: 68 variables

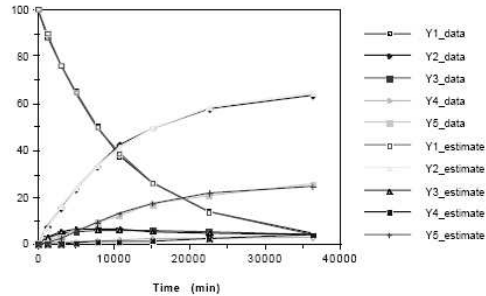
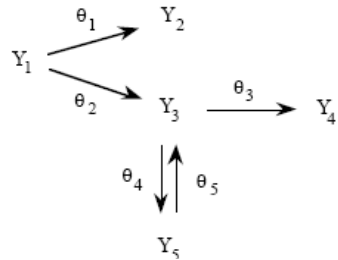
$(k_1, k_2, k_3)_0 = (6, 4, 1)$
 $(k_1, k_2, k_3)_* = (11.95, 7.99, 2.02)$
 $(k_1, k_2, k_3)_{true} = (12, 8, 2)$



Method	Obj.	Iters.	CPU (s, V3200)
BFGS SQP	8.23e-5	10	4.31
Gauss-Newton	8.23e-5	4	2.25
Hybrid SQP	8.23e-5	4	2.31

26

Small Residual Example: α -Pinene Kinetics



number of ODEs: 5
 number of parameters: 5
 discretized ODEs: 245 variables

Method	Iters.	CPU (s, V3200)
BFGS SQP	37	90.8
MINOS	21	64.8
Gauss-Newton	6	23.5
Hybrid SQP	6	23.6

27

Further Results (Tjoa and Biegler, 1990)

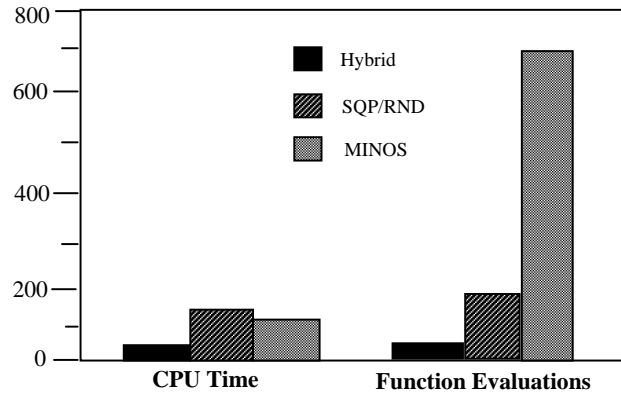
Number of Iterations (Function Evaluations)

Problem	GREG	MINOS	SQP based methods		
			BFGS	GN	Hybrid
1	6 (12)	8 (41)	8 (10)	3 (3)	3 (3)
2	14 (34)	8 (72)	14 (17)	3 (3)	3 (3)
3	6 (12)	8 (55)	30 (31)	3 (3)	3 (3)
4	7 (20)	13 (95)	14 (15)	6 (6)	6 (6)
5	7 (19)	15 (236)	fail	13 (16)	13 (16)
6	5 (10)*	7 (72)	13 (21)	5 (5)	5 (5)
7	8 (16)	10 (150)	24 (30)	5 (5)	6 (6)
8	fail	21 (271)	37 (61)	6 (6)	6 (6)
9	26 (95)*	33 (320)	19 (24)	fail	18 (18)

*: using least squares objective function

28

Comparison of Hybrid vs. General Purpose NLP Codes (Tjoa)



- Summary on 10 parameter estimation (kinetics) problems
- Few parameters, degrees of freedom
- Hybrid method for Hessian structure: (Fletcher and Xu) quasi-Newton method.

29

Further Comparison – Constrained Trust Region Method (Arora, 2003)

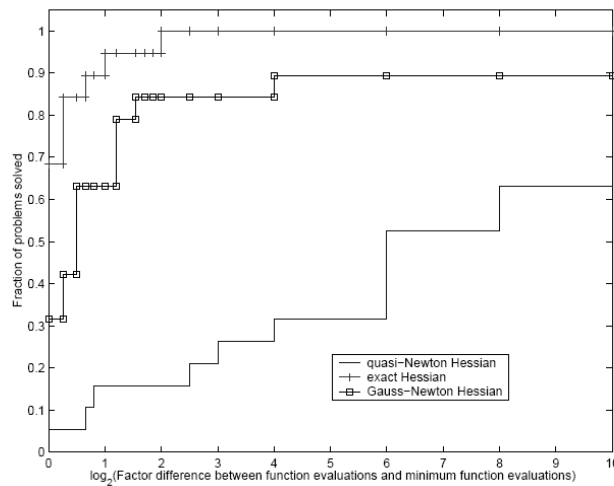


Figure 4.8: Parameter Estimation Problems

30

Statistical Inference of Estimated Parameters

Covariance of Optimal Parameters

Given an error distribution for the data (assumed Gaussian with covariance, V_z)
How does this affect the accuracy of the estimated parameters?

Recall: $V_z = E(\delta z \delta z^T)$ and $V_\theta = E(\delta \theta^* \delta \theta^{*T})$

How does θ change with data z ?

$$\text{i.e., } \frac{\partial \Phi(\theta^*, z)}{\partial \theta} = 0, \quad \frac{\partial \Phi(\theta^* + \delta \theta, z + \delta z)}{\partial \theta} = 0$$

$$\text{Approximate by: } \frac{\partial \Phi^2(\theta^*, z)}{\partial \theta \partial z} \delta z + \frac{\partial \Phi^2(\theta^*, z)}{\partial \theta^2} \delta \theta = 0$$

$$\text{and } \delta \theta = \left[\frac{\partial \Phi^2(\theta^*, z)}{\partial \theta^2} \right]^{-1} \frac{\partial \Phi^2(\theta^*, z)}{\partial \theta \partial z} \delta z, \text{ so we have:}$$

$$V_\theta = E(\delta \theta \delta \theta^T) = \left[\frac{\partial \Phi^2(\theta^*, z)}{\partial \theta^2} \right]^{-1} \frac{\partial \Phi^2(\theta^*, z)}{\partial \theta \partial z} E(\delta z \delta z^T) \frac{\partial \Phi^2(\theta^*, z)}{\partial \theta \partial z}^T \left[\frac{\partial \Phi^2(\theta^*, z)}{\partial \theta^2} \right]^{-T}$$

$$\Rightarrow V_\theta = \left[\frac{\partial \Phi^2(\theta^*, z)}{\partial \theta^2} \right]^{-1} \frac{\partial \Phi^2(\theta^*, z)}{\partial \theta \partial z} V_z \frac{\partial \Phi^2(\theta^*, z)}{\partial \theta \partial z}^T \left[\frac{\partial \Phi^2(\theta^*, z)}{\partial \theta^2} \right]^{-T}$$

31

Special Cases for Covariance

1. If the objective function has covariance independent of u , i.e.:

$$\Phi = \sum_{u=1}^n (z_u - f_u(\theta))^T V_z^{-1} (z_u - f_u(\theta)) \text{ then we have: } V_\theta = \sum_{u=1}^n (J_u^T V_z^{-1} J_u)^{-1}$$

2. If z_u is a scalar, then V_z is σ^2 and V_θ is given by: $\sigma^2 \sum_{u=1}^n \left(\frac{\partial f_u}{\partial \theta} \frac{\partial f_u}{\partial \theta}^T \right)^{-1}$

3. If the covariance for z is unknown, then estimate

$$\text{from moment matrix } V_z = M(\theta)/n \text{ and then: } V_\theta = \sum_{u=1}^n (J_u^T V_z^{-1} J_u)^{-1}$$

4. For general likelihood functions, $V_\theta = -(\nabla_{\theta\theta}(\log L(\theta^*)))^{-1}$ is asymptotically correct as $n \rightarrow \infty$.

32

Elliptical Confidence Regions

Single parameter

For a given interval let $\gamma = \Pr(a \leq \theta^* \leq b) = \int_a^b p(\theta^*|\theta_{\text{true}}) d\theta$

and for a single parameter this becomes: $|\theta_{\text{true}} - \theta^*| \leq \zeta$ where ζ is the confidence level for γ with σ_θ calculated with $n \rightarrow \infty$ and

$$\sigma_\theta^2 = \sigma^2 \sum_{u=1}^n \left(\frac{\partial f_u}{\partial \theta} \frac{\partial f_u}{\partial \theta}^T \right)^{-1}$$

Otherwise, with a small sample size: $|\theta_{\text{true}} - \theta^*| \leq t s_\theta$

$$s = \frac{1}{n-1} \left[\sum_{u=1}^n (z_u \cdot f_u(\theta))^2 \right]^{-1} \quad s_\theta^2 = s^2 \sum_{u=1}^n \left(\frac{\partial f_u}{\partial \theta} \frac{\partial f_u}{\partial \theta}^T \right)^{-1}$$

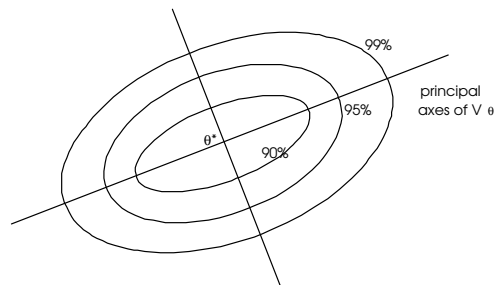
Multiple parameters

Map out a region $S(\theta)$ so that $\gamma = \Pr(\theta_{\text{true}} \in S(\theta))$

This can be done using the principal directions of V_θ which leads to:

$$\gamma = \Pr((\theta_{\text{true}} - \theta^*)^T V_\theta^{-1} (\theta_{\text{true}} - \theta^*))$$

33



For normal, unbiased distributions, linear models and a known V_θ , this probability follows a χ^2 distribution so that the region can be defined by:

$$(\theta_{\text{true}} - \theta^*)^T V_\theta^{-1} (\theta_{\text{true}} - \theta^*) \leq c(\gamma)$$

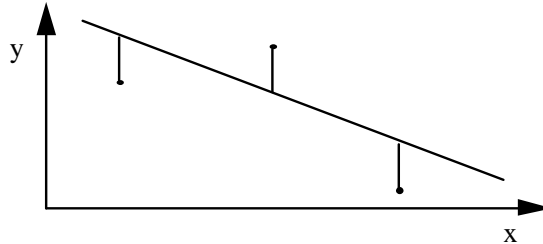
$c(\gamma)$ is χ^2 value for γ level of confidence with n_θ degrees of freedom.

- For a scalar z , the χ^2 test simplifies to an F-test for determination of $c(\gamma)$.
- Elliptical confidence regions are correct if the model is linear or for small levels of confidence, γ . Otherwise, confidence regions can deviate greatly from ellipses.
- Elliptical confidence regions are most commonly used. Nonlinear confidence regions much more expensive to calculate.

34

Errors in Variables Models (EVM)

Conventional model: $y = f(x, \theta)$ or $h(x, y, \theta) = 0$



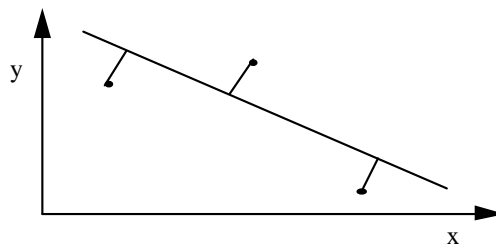
- No. of independent variables (x) = s_1
- No. of dependent variables (y) = s_2
- No. of parameters (θ) = p
- No. of constraints = m
- No. of data sets = r
- \Rightarrow No error in the "independent" variables, x
- Number of degrees of freedom for optimization = p
- Minimize vertical distance, e.g. $(y_u - z_u)^2$

35

Errors in Variables Models (EVM) - 2

Implicit Model: $f(x, y, \theta) = 0$

Both x and y have inherent measurement errors
 \Rightarrow under-determined system
 (e.g., pressure vs. temperature data)



- Number of degrees of freedom for optimization
 $= p + (s-m)r$
- \Rightarrow NLP size grows linearly with the number of data sets
- Minimize a *nonvertical* distance

36

EVM Problem Statement

$$\text{Min } \Phi = 1/2 \sum_{u=1}^r e_u^T W e_u$$

$$\text{s.t. } f_u(w_u, \theta) = 0 \\ a \leq \theta \leq b$$

$$\text{where } e_u = w_u - z_u, w_u^T = [x_u^T y_u^T]$$

Formulation:

- Least squares nonlinear constraints with many degrees of freedom

Current Approaches

- Nonlinear Programming Strategies – Expensive if sparsity not exploited
- Linearized Least Squares - Not robust, global convergence not enforced

$$\text{Min } \Phi = \sum_{u=1}^r \phi_u = 1/2 \sum_{u=1}^r e_u^T W e_u$$

$$\text{s.t. } h_u(w_u, \zeta_u, \theta) = 0 \\ a \leq \theta \leq b$$

$$\text{where } h_u^T = [f_u^T g_u^T]$$

Introduce

- New variables: ζ_u
- New constraints: $g_u = \zeta_u - \theta = 0$.

Solution Strategy:

- Apply decoupling strategy to SQP method
- Decompose for each data set
- Computational cost linear in # of data sets
- Can solve in parallel

37

KKT Conditions for EVM

Structure of the KKT matrix after decoupling:

= RHS($\Delta\theta$)

$$\begin{bmatrix} B_u & \nabla h_u \\ \nabla h_u^T & 0 \end{bmatrix} \begin{bmatrix} d_u \\ v_u \end{bmatrix} = - \begin{bmatrix} \nabla \phi_u \\ h_u \end{bmatrix} \quad \nabla h_u = \begin{bmatrix} \nabla_w f_u & 0 \\ \nabla_\zeta f_u & I \end{bmatrix} \quad d = \begin{bmatrix} \Delta w_u \\ \Delta \zeta_u \end{bmatrix} \quad v_u = \begin{bmatrix} \lambda_u \\ \gamma_u \end{bmatrix}$$

Once θ is fixed, can solve for remaining variables independently.

Can exploit each KKT system further.

$$\text{Structure of Hessian (Gauss-Newton): } B_u = \begin{bmatrix} W_u & 0 \\ 0 & 0 \end{bmatrix}$$

38

EVM Decomposition

$$\begin{bmatrix} B_u & \nabla h_u \\ \nabla h_u^T & 0 \end{bmatrix} \begin{bmatrix} d_u \\ v_u \end{bmatrix} = - \begin{bmatrix} \nabla \phi_u \\ h_u \end{bmatrix}$$

Apply Range and Null space Decomposition

- Subspace search directions: $\nabla h_u^T Z_u = 0, Y_u^T Z_u = 0$
- Define: $d_u = Z_u p_{Z,u} + Y_u p_{Y,u}$
 Null space step: $p_{Z,u} = - (Z_u^T B_u G N Z_u)^{-1} Z_u^T \nabla \phi_u$
 Range space step: $p_{Y,u} = - (\nabla h_u^T Y_u)^{-1} h_u$
- Note: $p_{Y,u}$ is dependent on $\Delta\theta$, $p_{Z,u}$ is not!

Reconstruct QP Problem is space of $\Delta\theta$

- Sum up contributions from all data sets

$$\text{Min}_{\Delta\theta} \left(\sum_{u=1}^r \alpha_u \right)^T \Delta\theta + 1/2 \Delta\theta^T \left(\sum_{u=1}^r \beta_u \right) \Delta\theta$$

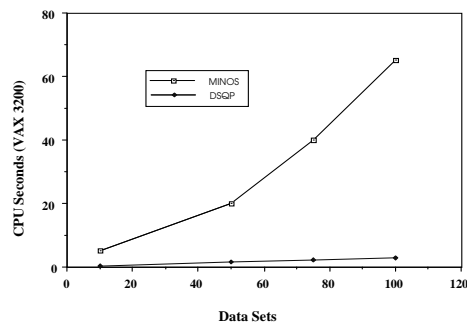
$$\text{s.t. } a - \theta^k \leq \Delta\theta \leq b - \theta^k$$

where α_u and β_u are derived from ∇h_u

39

EVM Examples

Example	Model	Data sets	Number of d.o.f
1	$y = \theta_1 + \theta_2 x$	10	12
2	$y = \theta_1 + \theta_2 x + \theta_3 x^2 + \theta_4 x^3$	10	14
3	$\exp(-\theta_1 t) \exp(-\theta_2/T) - y = 0$	15	32
4	$x_3 = \frac{\theta_1 \theta_2^2 \theta_3 x_1 x_2^2}{(1 + \theta_1 x_1 + \theta_2 x_2)^3}$	28	59
5, Case 1	$x_2 = \theta_1 + \frac{1}{x_1 - \theta_2}$	25	27
5, Case 2		50	52
5, Case 3		75	77
5, Case 4		100	102



40

Chemical Engineering

Full-space NLP Formulation for Parameter Estimation

Original Formulation

$$\begin{aligned} \min_{x \in \mathcal{R}^n} & f(x) \\ \text{s.t.} & c(x) = 0 \\ & x \geq 0 \end{aligned}$$

Can generalize for

$$a \leq x \leq b$$

Barrier Approach

$$\begin{aligned} \min_{x \in \mathcal{R}^n} & \varphi_\mu(x) = f(x) - \mu \sum_{i=1}^n \ln x_i \\ \text{s.t.} & c(x) = 0 \end{aligned}$$

Fiacco and McCormick (1968)

⇒ As $\mu \rightarrow 0$, $x^*(\mu) \rightarrow x^*$

41

Chemical Engineering

Solution of the Barrier Problem - IPOPT

⇒ Newton Directions (KKT System)

$$\begin{aligned} \nabla f(x) + A(x)\lambda - v &= 0 \\ XVe - \mu e &= 0 \\ c(x) &= 0 \end{aligned}$$

⇒ Reducing the System

$$d_v = \mu X^{-1} e - v - X^{-1} V d_x$$

$$\begin{bmatrix} Q + \Sigma & A \\ A^T & 0 \end{bmatrix} \begin{bmatrix} d_x \\ \lambda^+ \end{bmatrix} = - \begin{bmatrix} \nabla \varphi_\mu \\ c \end{bmatrix} \quad \Sigma = X^{-1} V$$

What are the Benefits for Parameter Estimation?

42

Chemical Engineering

Large-Scale Parameter Estimation

Material & Energy $\left\{ \begin{array}{l} F_j \left[\frac{dy_j(z)}{dz}, y_j(z), w_j(z), z, \pi_j, \Pi \right] = 0 \\ \text{Physical Properties} \left\{ \begin{array}{l} G_j \left[y_j(z), w_j(z), z, \pi_j, \Pi \right] = 0 \\ \text{Zone Transitions} \left\{ \begin{array}{l} y_j(0) = \phi(y_{j-1}(z_{L_{j-1}}), F_{fj}) \end{array} \right. \end{array} \right. \right. \right. j \in \{1..NZ\}$

500 ODEs
1000 AEs

× Stiffness + Highly Nonlinear + Parametric Sensitivity + Algebraic Coupling

CarnegieMellon 43

Chemical Engineering

Large-Scale Parameter Estimation

□ Complex Kinetic Mechanisms

Initiator decomposition
 $I_i \xrightarrow{k_{d_i}} 2R \quad i = 1, N_I$

Chain Initiation
 $R + M_1 \xrightarrow{k_{i1}} P_{1,0}$
 $R + M_2 \xrightarrow{k_{i2}} Q_{0,1}$

Chain Propagation
 $P_{r,s} + M_1 \xrightarrow{k_{p11}} P_{r+1,s}$
 $P_{r,s} + M_2 \xrightarrow{k_{p12}} Q_{r,s+1}$
 $Q_{r,s} + M_1 \xrightarrow{k_{p21}} P_{r+1,s}$
 $Q_{r,s} + M_2 \xrightarrow{k_{p22}} Q_{r,s+1}$

Chain Transfer to Monomer
 $P_{r,s} + M_1 \xrightarrow{k_{fm11}} P_{1,0} + M_{r,s}$
 $P_{r,s} + M_2 \xrightarrow{k_{fm12}} Q_{0,1} + M_{r,s}$
 $Q_{r,s} + M_1 \xrightarrow{k_{fm21}} P_{1,0} + M_{r,s}$
 $Q_{r,s} + M_2 \xrightarrow{k_{fm22}} Q_{0,1} + M_{r,s}$

Chain Transfer to Solvent
 $P_{r,s} + S_i \xrightarrow{k_{s1i}} P_{1,0} + M_{r,s}$
 $Q_{r,s} + S_i \xrightarrow{k_{s2i}} Q_{0,1} + M_{r,s}$

Chain Transfer to Polymer
 $P_{r,s} + M_{x,y} \xrightarrow{k_{fp11}} P_{x,y} + M_{r,s}$
 $P_{r,s} + M_{x,y} \xrightarrow{k_{fp12}} Q_{x,y} + M_{r,s}$
 $Q_{r,s} + M_{x,y} \xrightarrow{k_{fp21}} P_{x,y} + M_{r,s}$
 $Q_{r,s} + M_{x,y} \xrightarrow{k_{fp22}} Q_{x,y} + M_{r,s}$

Termination by Combination
 $P_{r,s} + P_{x,y} \xrightarrow{k_{tc11}} M_{r+x,s+y}$
 $P_{r,s} + Q_{x,y} \xrightarrow{k_{tc12}} M_{r+x,s+y}$
 $Q_{r,s} + Q_{x,y} \xrightarrow{k_{tc22}} M_{r+x,s+y}$

Termination by Disproportionation
 $P_{r,s} + P_{x,y} \xrightarrow{k_{td11}} M_{r,s} + M_{x,y}$
 $P_{r,s} + Q_{x,y} \xrightarrow{k_{td12}} M_{r,s} + M_{x,y}$
 $Q_{r,s} + Q_{x,y} \xrightarrow{k_{td22}} M_{r,s} + M_{x,y}$

Backbiting
 $P_{r,s} \xrightarrow{k_{b1}} P_{r,s} \text{ or } Q_{r,s}$
 $P_{r,s} \xrightarrow{k_{b2}} Q_{r,s} \text{ or } P_{r,s}$

β-scission
 $P_{r,s} \xrightarrow{k_{\beta 1}} M_{r,s}^= + P_{1,0}$
 $P_{r,s} \xrightarrow{k_{\beta 2}} M_{r,s}^= + Q_{0,1}$

$k = k_0 \exp\left(-\frac{E_a + P E_v}{RT}\right)$

~ 35 Elementary Reactions
~ 100 Kinetic Parameters

CarnegieMellon 44

Large-Scale Parameter Estimation

- Parameter Estimation for Industrial Applications
 - Use Rigorous Model to Match Plant Data Directly
 - Start with Standard Least-Squares Formulation

$$\begin{aligned}
 \min_{\Pi, \pi_{k,j}} & \sum_{k=1}^{NS} \sum_{j=1}^{NZ} \sum_{i=1}^{NM(j)} (y_{k,j}(z_i) - y_{k,j,i}^M)^T V_y^{-1} (y_{k,j}(z_i) - y_{k,j,i}^M) \\
 & + \sum_{k=1}^{NS} (w_{k,NZ} - w_{k,NZ}^M)^T V_w^{-1} (w_{k,NZ} - w_{k,NZ}^M) \quad \left. \vphantom{\min} \right\} \text{Least-Squares} \\
 \text{s.t.} & \\
 & F'_{k,j} \left[\frac{dy_{k,j}(z)}{dz}, y_{k,j}(z), w_{k,j}(z), z, \pi_{k,j}, \Pi \right] = 0 \\
 & G_{k,j} [y_{k,j}(z), w_{k,j}(z), z, \pi_{k,j}, \Pi] = 0 \quad \left. \vphantom{\min} \right\} \text{Rigorous Reactor Model} \\
 & y_{k,j}(0) = \phi(y_{k,j-1}(z_{L_{k,j-1}}), F_{f_{k,j}}) \\
 & j \in \{1..NZ\}, \quad k \in \{1..NS\}
 \end{aligned}$$

- Special Case of Multi-Stage Dynamic Optimization Problem
 - Solve using Simultaneous Collocation-Based Approach

1 data set		6 data sets
500 ODEs	x 6	3000 ODEs
1000 AEs	→	6000 AEs _s

Large-Scale Parameter Estimation

- Multi-Zone Tubular Reactor – Quasi Steady-State
 - Data Sets: Operating Conditions and Properties for Different Grades
 - Match: Temperature Profiles and Product Properties
 - On-line Adjusting Parameters → Track Evolution of Disturbances
 - Kinetic Parameters → Development and Discrimination among Rigorous Models

□ Results

□ Single Data Set (On-line Adjusting Parameters)

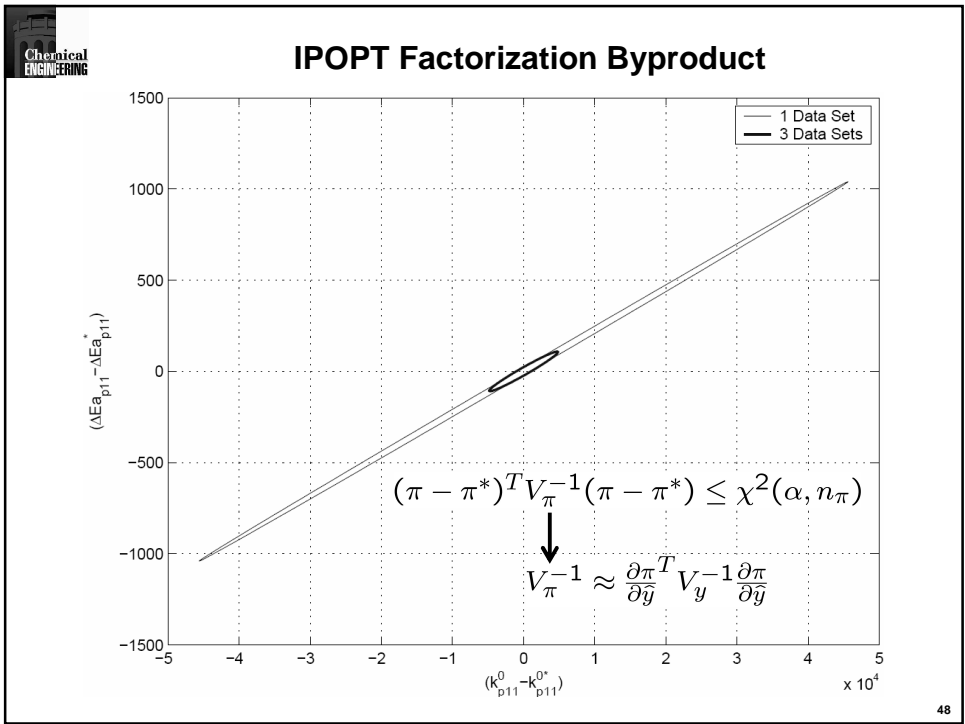
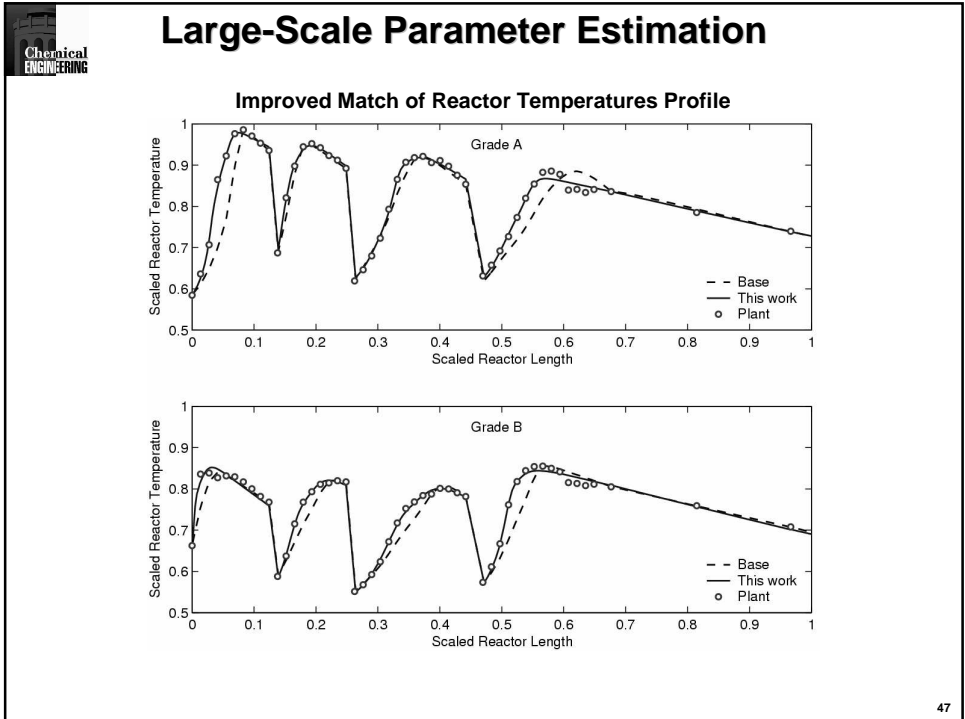
Grade	Constraints	Parameters	LB	UB	Iterations	CPU _s	NZJ	NZH
A	11955	32	374	361	11	17.03	166425	87954
B	11283	32	374	361	8	10.06	138666	76890

□ Multiple Data Sets (On-line Adjusting Parameters + Kinetics)

Data Sets	Constraints	DOF	LB	UB	Iterations	CPU _s	NZJ	NZH
3	33900	121	1246	1207	68	451.51	520275	552738
6	68421	217	2467	2389	58	900.21	1058412	1119258

**Bottleneck (Memory Requirements)
Factorization Step**

$$\begin{bmatrix} W(x_k, \lambda_k) & A(x_k) & -I \\ A(x_k)^T & 0 & 0 \\ V_k & 0 & X_k \end{bmatrix} \begin{bmatrix} \Delta x \\ \Delta \lambda \\ \Delta \nu \end{bmatrix} = - \begin{bmatrix} \nabla f(x_k) + A(x_k) \lambda_k - \nu_k \\ c(x_k) \\ X_k V_k e - \mu e \end{bmatrix}$$



Advanced Regression Methods

Errors-In-Variables (EVM)

- Standard Least Squares - Errors in Output Variables - *Biased* Parameters
- EVM - Errors in Output *AND* Input Variables - *Unbiased* Parameters

Inputs

Inputs

Outputs

Carnegie Mellon 49

Advanced Regression Methods

Errors-In-Variables (EVM)

EVM Drawback - Degrees of Freedom $DOF = \Pi + \sum_{k=1}^{NS} \sum_{j=1}^{NZ} \pi_{k,j} + \sum_{k=1}^{NS} \sum_{j=1}^{NZ} \sum_{i=1}^{NM_u(j)} u_{k,j,i}$

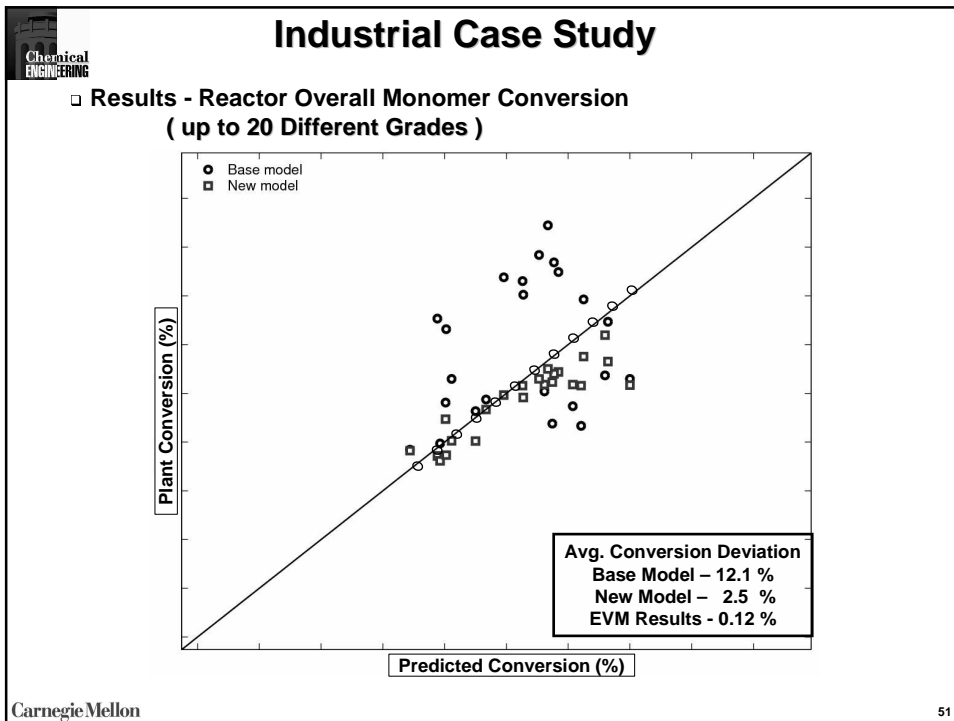
Formulation is Straightforward

$$\min_{\Pi, \pi_{k,j}, u_{k,j}} \sum_{k=1}^{NS} \sum_{j=1}^{NZ} \sum_{i=1}^{NM(j)} (y_{k,j}(z_i) - y_{k,j,i}^M)^T \mathbf{V}_y^{-1} (y_{k,j}(z_i) - y_{k,j,i}^M) + \sum_{k=1}^{NS} (w_{k,NZ} - w_{k,NZ}^M)^T \mathbf{V}_w^{-1} (w_{k,NZ} - w_{k,NZ}^M) + \sum_{k=1}^{NS} \sum_{j=1}^{NZ} (u_{k,j} - u_{k,j}^M)^T \mathbf{V}_u^{-1} (u_{k,j} - u_{k,j}^M)$$

EVM vs. Standard Least Squares

Data Sets	Constraints	DOF	LB	UB	Iterations	CPUs	NZJ	NZH
6 (EVM)	68627	529	2653	2575	71	1010.74	1059512	1119780
6 (SLS)	68421	217	2467	2389	58	900.21	1058412	1119258

50



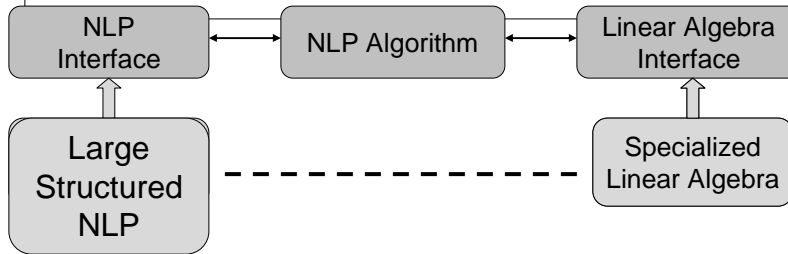
Chemical Engineering

Next Generation IPOPT

- IPOPT 3.2, Fall, 2006, CPL (www.coin-or.org)
 - Based on Fortran version
 - Object-oriented, NLP Solver
 - Primal-Dual Interior Point method
 - Full space - exact Hessian information
 - Monotone/Adaptive μ update
 - Filter line search strategy
 - Flexible algorithm structure
 - Interfacing to other linear solvers
 - Modeling Environments - AMPL, AIMMS, MATLAB...
- *Ideal for Internal Decomposition*
 - Consistent linear system structure at every iteration
 - Separation of algorithm and specialized linear algebra



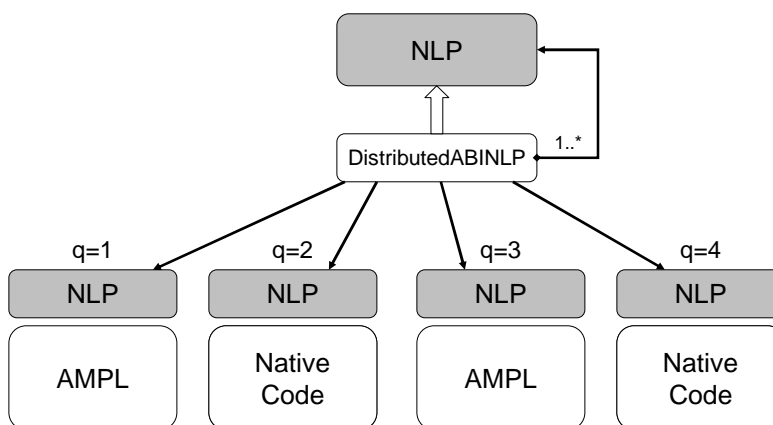
High Level IPOPT Design



- Provide structured NLP elements
 - Hessian, jacobian, gradients, residuals, variables,
- Vector operations
 - BLAS, norms, dot products, axpy, element- wise, max, min, etc.
- Matrix-Vector operations
 - Mv , $M^T v$
- Solution of Linear KKT system



Modeling the Block NLP



54

Parameter Estimation in Parallel Architectures

Exploit Structure of KKT Matrix – Laird, Biegler 2006

$$\min_{\Pi, \pi_{k,j}} \sum_{k=1}^{NS} \sum_{j=1}^{NZ} \sum_{i=1}^{NM(j)} (y_{k,j}(z_i) - y_{k,j,i}^M)^T V_y^{-1} (y_{k,j}(z_i) - y_{k,j,i}^M) + \sum_{k=1}^{NS} (w_{k,NZ} - w_{k,NZ}^M)^T V_w^{-1} (w_{k,NZ} - w_{k,NZ}^M)$$

s.t.

$$F_{k,j} \begin{bmatrix} \frac{dy_{k,j}(z)}{dz}, y_{k,j}(z), w_{k,j}(z), z, \pi_{k,j}, \Pi \end{bmatrix} = 0$$

$$G_{k,j} \begin{bmatrix} y_{k,j}(z), w_{k,j}(z), z, \pi_{k,j}, \Pi \end{bmatrix} = 0$$

$$y_{k,j}(0) = \phi(y_{k,j-1}(z_{L_{k,j-1}}), F_{f_{k,j}})$$

$$j \in \{1..NZ\}, k \in \{1..NS\}$$

min $f(x)$

s.t. $c(x) = 0$

$x \geq 0$

$$\begin{bmatrix} H_k & A_k & -I \\ A_k^T & 0 & 0 \\ V_k & 0 & X_k \end{bmatrix} \begin{bmatrix} \Delta x \\ \Delta \lambda \\ \Delta \nu \end{bmatrix} = - \begin{bmatrix} r_x \\ r_\lambda \\ X_k V_k e - \mu_\ell e \end{bmatrix}$$

Direct Factorization MA27

Memory Bottlenecks

Factorization Time Scales Superlinearly with Data sets

min $\sum_{k=1}^{NS} f_k(x_k)$

s.t. $c_k(x_k, \Pi) = 0$

$x_k, \Pi \geq 0$

$$\begin{bmatrix} K_1 & & & Q_1 \\ & K_2 & & Q_2 \\ & & \dots & \vdots \\ & & & K_N \\ Q_1^T & Q_2^T & \dots & Q_N^T \end{bmatrix} \begin{bmatrix} \Delta s_1 \\ \Delta s_2 \\ \vdots \\ \Delta s_N \\ \Delta \Pi \end{bmatrix} = - \begin{bmatrix} r_1 \\ r_2 \\ \vdots \\ r_N \\ r_\Pi \end{bmatrix}$$


Block-bordered Diagonal Structure

Coarse-Grained Parallelization

Carnegie Mellon 55

Parameter Estimation in Parallel Architectures

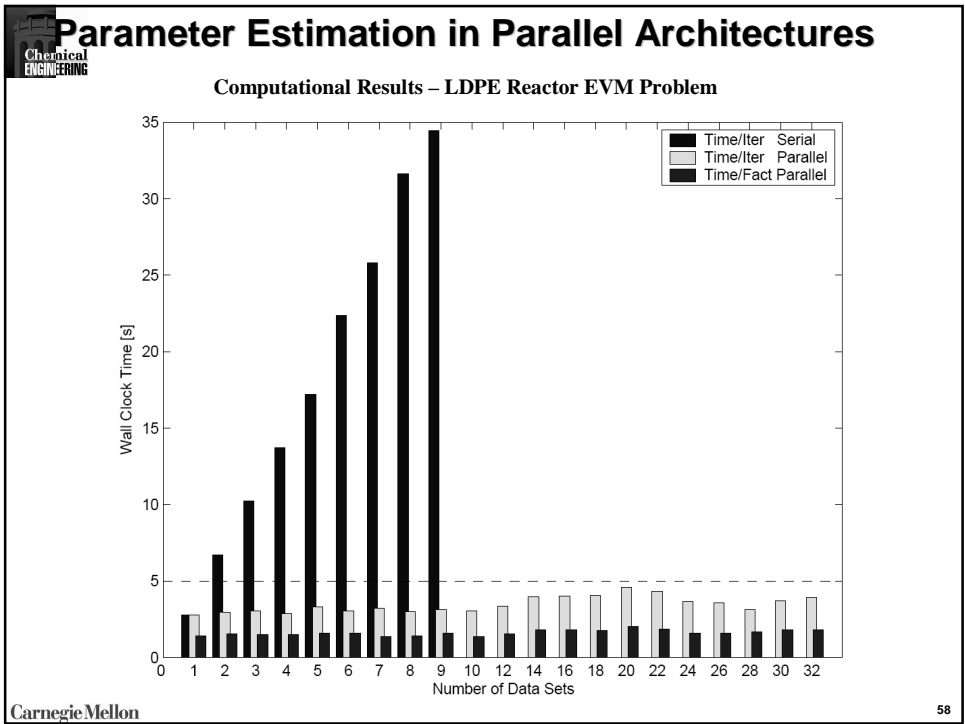
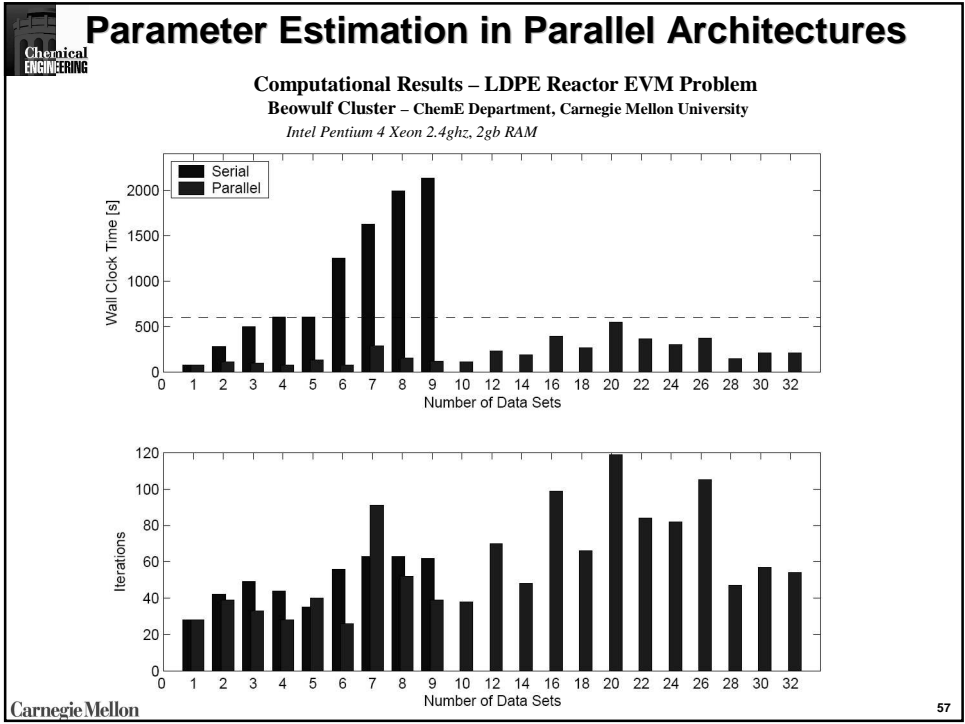
Sch



$$\begin{bmatrix} \Delta s_1 \\ \Delta s_2 \\ \vdots \\ \Delta s_N \\ \Delta \Pi \end{bmatrix} = - \begin{bmatrix} r_1 \\ r_2 \\ \vdots \\ r_N \\ r_\Pi \end{bmatrix}$$

- 1) **Parallelizable**
Sparse Factorization K_i in each Block
Inertia Correction
- 2) **Non-parallelizable**
Dense Factorization of $S \rightarrow \dim\{\Pi\}$
- 3) **Parallelizable**
Final Backsolve

Carnegie Mellon 56



Conclusions – Parameter Estimation

- Trust Region (Levenberg-Marquardt) methods: standard for unconstrained problems - MINPACK, NaG, NL2SOL, Harwell
- Constrained problem formulations - more model flexibility
- SQP codes adapted to exploit least squares structure, faster methods
- EVM problems - expensive for conventional optimization codes
 - Many degrees of freedom for optimization
 - Decomposition of KKT conditions required
 - ODRPACK (netlib) developed for $y_u = f_u(\theta)$
- Large-scale SQP methods developed for:
 - Parameter estimation
 - EVM methods
 - Data Reconciliation
- IPOPT has useful characteristics for large-scale parameter estimation

59

Optimization Algorithms for Data Reconciliation

Introduction to Data Reconciliation

M-Estimators

- Bayesian Forms
- Fair Function
- Redescending Estimator

Akaike Information Criterion

Mixed Integer Formulations

Static Examples

Dynamic Examples

Conclusions

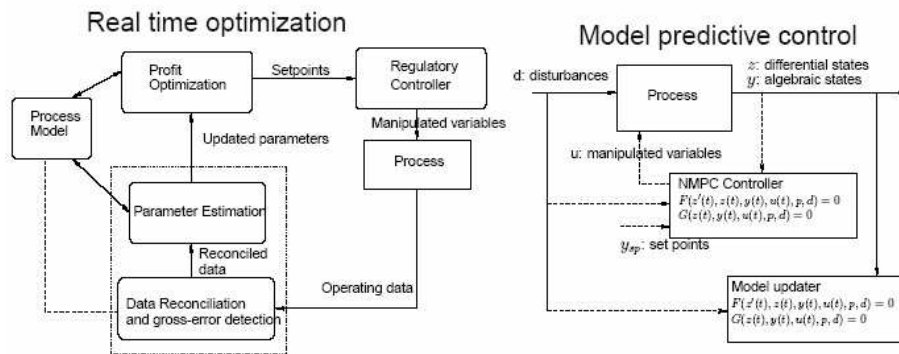
60

Introduction

- Data contaminated by errors
 - Random errors:
 - High frequency noise
 - Typically $\sim \mathcal{N}(0, \Sigma^2)$
 - Gross errors and outliers
 - Systematic occurrences: biases, drifts
 - Not modeled by distributions
- Contaminated data can bias estimates
- Difficult to converge NLP when data noisy
- Inputs may be contaminated

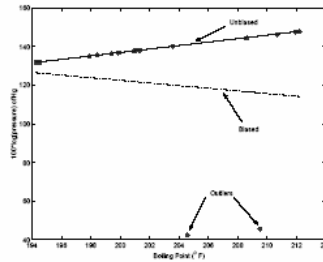
Data Reconciliation Framework

- Industrial applications of data reconciliation and parameter estimation (DRPE)



Effects of Gross Error in Regression

Example: Boiling point of H₂O vs. atm. pressure with 2 outliers (Weisberg, 1985)



Sources of gross errors: Broken gauges, process leaks, improperly used measurement devices, operator errors

Identify and negate effects of gross errors

63

Data Reconciliation – Literature Review

- Sequential χ^2 tests:
 - Crowe et al. (1983), Madron (1985), Kao et al. (1992)
- Combinatorial
 - Serth and Heenan (1986), Narasimhan and Mah (1987): Serial Combination and Generalized Likelihood Ratio tests
 - Crowe (1989): Max. power test
 - Rollins et al. (1996): Linear Combination Test based on χ^2 tests
 - Bagajewicz et al. (1998): Graph theory based, linear ODEs
- Simultaneous
 - Tjoa and Biegler (1991): Contaminated normal - difficult to tune, cannot detect biases
 - Albuquerque and Biegler (1996): Fair function - easy to use, no explicit outlier identification

64

Treatment of Outliers: M-Estimators

M-Estimator: modeling the influence of residual outliers through modification of maximum likelihood (ML) functions

Bayesian Statistics and Bernoulli Trials

- Statistical definitions and ML function
- Allows statistical inference

Robust Statistics

- Modification without inferential aspects

Akaike Information Criterion (AIC)

- Based on ML extended to discrete parameters

65

Bayesian Approach: Bernoulli Trial

Assume separate probability distributions for random and gross errors:

$$P(\varepsilon_i | \theta, X_i) = R(\varepsilon_i | \theta, X_i)\Pi(R) + G(\varepsilon_i | \theta, X_i)\Pi(G)$$

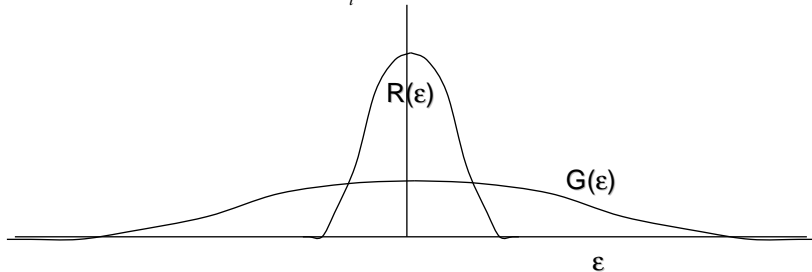
P, R, G – probability distributions

$\Pi(R), \Pi(G)$ – prior distributions

ε_i – measurement error

θ – parameters

X_i – measurements



66

Bayesian Approach: Bernoulli Trial

Combine both distributions into ML function:

$$P(\varepsilon_i | \theta, X_i) = R(\varepsilon_i | \theta, X_i)\Pi(R) + G(\varepsilon_i | \theta, X_i)\Pi(G)$$

$$P(\varepsilon_i | \theta, X_i) = R(\varepsilon_i | \theta, X_i)(1-p) + G(\varepsilon_i | \theta, X_i)p$$

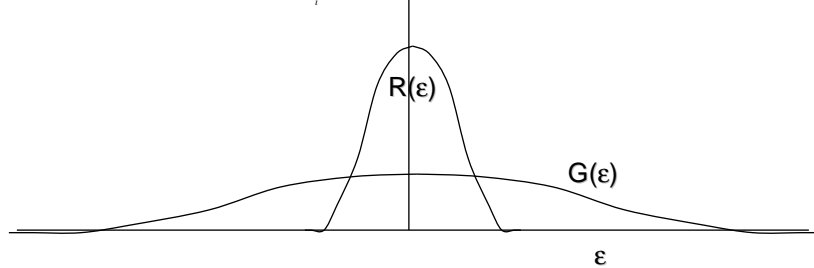
$$P(\varepsilon_i | \theta, X_i) = \exp(-\varepsilon_i^2 / \sigma_i^2)(1-p) / \sigma_i + \exp(-\varepsilon_i^2 / (b\sigma_i)^2)p / (b\sigma_i)$$

$$\max \sum_i \log(P(\varepsilon_i | \theta, X_i))$$

p - fraction gross errors

b - scale of distributions

σ_i^2 - var measurement i



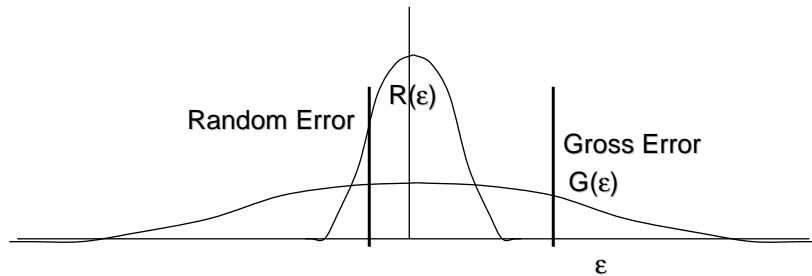
67

Bayesian Approach: Gross Error Test

Solve optimization problem and perform gross error tests at solution:

$$\max \sum_i \log(\exp(-\varepsilon_i^2 / \sigma_i^2)(1-p) + \exp(-\varepsilon_i^2 / (b\sigma_i)^2)p / b)$$

$$\text{If gross error: } \exp(-\varepsilon_i^2 / \sigma_i^2)(1-p) < \exp(-\varepsilon_i^2 / (b\sigma_i)^2)p / b$$



- + Statistical basis for determining gross errors
- Assumes gross errors follow proposed distribution
- Not robust to deviations from assumptions

68

Example (Pai and Fisher, 1988)

$$0.5(x_1)^2 - 0.7x_2 + x_3u_1 + (x_2)^2u_1u_2 + 2x_3(u_3)^2 - 255.8 = 0$$

$$x_1 - 2x_2 + 3x_1x_2 - 2x_2u_1 - x_2u_2u_3 + 111.2 = 0$$

$$x_3u_1 - x_1 + 3x_1x_3 - 2x_2u_1 - x_2u_2u_3 + 111.2 = 0$$

$$x_4 - x_1 - (x_3)^2 + u_2 + 3u_3 = 0$$

$$x_5 - 2x_3u_2u_3 = 0$$

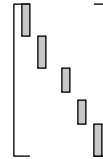
$$2x_1 + x_2x_3u_1 + u_2 - u_3 - 126.6 = 0$$

5 measured variables (x)

3 unmeasured variables (u)

b= 20, p=0.2, σ=0.1

**Introduce 20% gross error into
500 data sets generated randomly**



**100 right
0 wrong**



**87 right
0 wrong**



**99 right
1 wrong**

69

Robust Statistics for Gross Errors

- Common Assumption: Data with normally distributed random errors
⇒ likelihood function is *least squares estimator*
- likelihood of data corrupted with outliers: difficult to determine
- Estimators derived from fixed probability distribution: not justified
- Robust Estimators:
 - largely distribution independent
 - unbiased estimates when data from ideal distribution
 - Insensitive to deviations from ideality
 - Add less weight to outlying measurements ⇒ *protect* other measurements from being corrupted

70

Robust Statistics Properties

Find estimators that are insensitive to deviations in assumptions of noise distribution.

F – assumed distribution of data

G – actual distribution of data $\delta(F, G) < \varepsilon \Rightarrow \delta(T(F), T(G)) < \gamma$

T(-) – test statistic or distribution of estimator

Example:

S1: [2.0, 2.1, 2.2, 2.3, 2.4] Mean = 2.2, Median=2.2

S2: [2.0, 2.1, 2.2, 2.3, 24] Mean = 6.56, Median=2.2

→ Median is a robust statistic, Mean is not.

71

Huber Class of Robust Estimators (1981)

Fair Function (Rey, 1988)

$$\rho_{Fi} = C^2 \left[\frac{|\varepsilon_i|}{C} - \log \left\{ 1 + \frac{|\varepsilon_i|}{C} \right\} \right]$$

- Use $\sum_i \rho_{Fi}$ instead of least squares objective
- Convex function of normalized residual
- Small residuals – quadratic behavior
- Large residuals – linear behavior
- C – Tuning parameter based on Cramer-Rao bound (trade off efficiency (high C) with robustness (low C))

$$C = 0.21529 \left(\frac{\varphi - 0.63662}{1 - \varphi} \right) \text{ where } \varphi = \frac{\sigma^2}{V_{Fair}}$$

72

Hampel Class of Robust Estimators

Three Part Redescending Estimator (Hampel, 1974)

$$\rho_{H_i} = \begin{cases} \frac{1}{2}\varepsilon_i^2, & 0 \leq |\varepsilon_i| \leq a \\ a|\varepsilon_i| - \frac{a^2}{2}, & a < |\varepsilon_i| \leq b \\ ab - \frac{a^2}{2} + (c-b)\frac{a}{2} \left[1 - \left(\frac{c-|\varepsilon_i|}{c-b} \right)^2 \right], & b < |\varepsilon_i| \leq c \\ ab - \frac{a^2}{2} + (c-b)\frac{a}{2}, & |\varepsilon_i| > c, \end{cases}$$

Three parameters define regions: $c > b+2a$

Quadratic, linear and constant parts involving normalized residual

Nonsmooth – requires some smoothing approximations

Not clear how to tune a, b, c

73

M-Estimators in Robust Statistics

M-Estimators

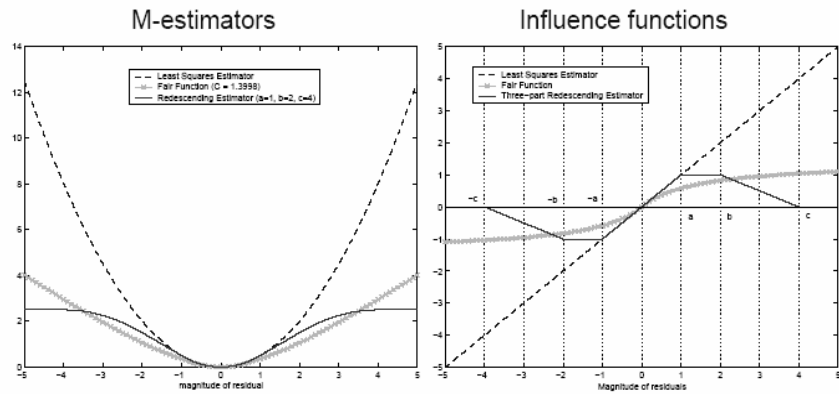
- Overall likelihood of N observations: $\mathcal{L} = \prod_{i=1}^N \ell(x_i|p)$
- M-estimator = $\rho^M = -\log \mathcal{L} = -\sum_{i=1}^N \log \ell(x_i|p)$

M-Estimators Used $\mathcal{IF} = \psi(\varepsilon_0) = \lim_{t \rightarrow 0} \frac{T[(1-t)f + t\delta(\varepsilon - \varepsilon_0)] - T[f]}{t}$

- Least Squares Estimator (Non robust)
 - Severely biased by outliers
- Fair Function (Robust)
 - \mathcal{IF} is bounded as $\varepsilon \rightarrow \infty$
- Three part Redescending Estimator (Robust)
 - $\mathcal{IF} \rightarrow 0$ as $\varepsilon \rightarrow \infty$

74

Properties of M-Estimators



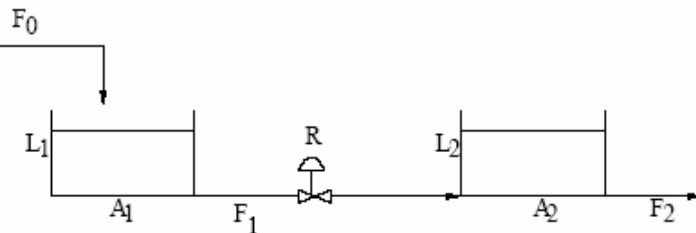
75

Properties of M-Estimators

	Bayesian	Robust
Distributional Assumptions?	Yes	No
Sensitive to Assump. Deviations?	Yes	No
Statistical Inference at Solution?	Yes	No
Incorporate prior knowledge?	Yes	No
-For robust case, apply exploratory statistics (e.g., boxplots) at solution		
- Need data redundancy to identify gross errors		

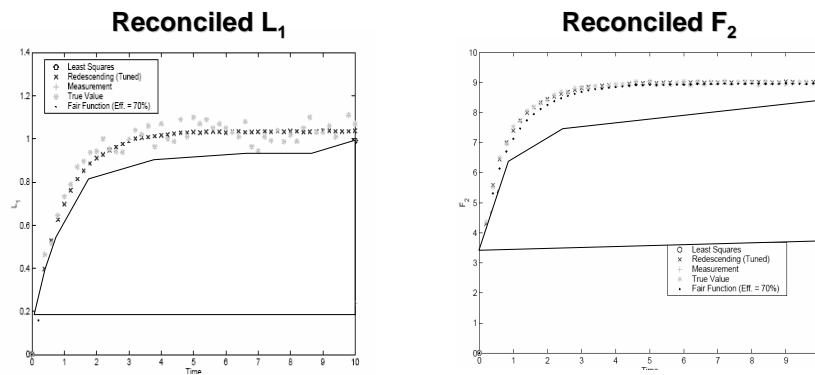
76

Dynamic Data Reconciliation



- Measurements: Flows F_0 , F_1 , F_2 & heights L_1 and L_2
- DAE system
- Parameters: Tank areas: $1/A_1$ & $1/A_2$

Case 1: Comparison of Results



- Discretize DAEs and solve as large scale NLP problem with appropriate objective function
- Add random outlier noise to flow and level data

Tank Example – Case 1

- Outliers drawn from broad random distribution (distributional assumption satisfied)
- Data reconciliation and parameter estimates done well by
- Bayesian approach and M-estimators

<u>Parameters</u>	<u>1/A₁</u>	<u>1/A₂</u>
True Values	0.5	0.5
Least Squares	0.698	0.503
Bayesian	0.490	0.50
Fair Function	0.501	0.501
Redesc. Tuned	0.50	0.50

79

Tank Example – Case 2

- Outliers systematic (distributional assumption violated); measurements for L₁ and F₂ stuck
- Data reconciliation and parameter estimation poor with Bayesian approach
- Data reconciliation does and parameter estimates done better by
- M-estimators

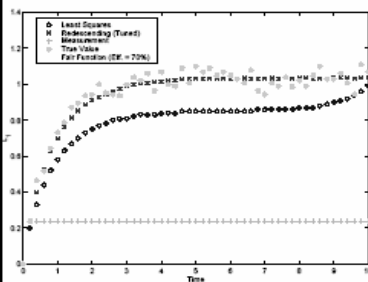
<u>Parameters</u>	<u>1/A₁</u>	<u>1/A₂</u>
True Values	0.5	0.5
Least Squares	0.25*	0.55
Bayesian	0.25	0.25
Fair Function	0.25	0.439
Redesc. Tuned	0.499	0.500

*lower bound

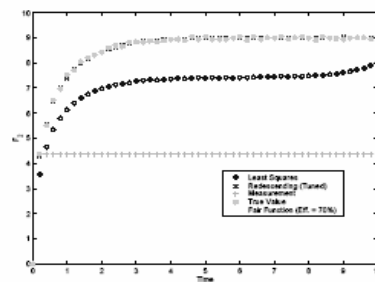
80

Case 2: Comparison of Results

Reconciled L_1



Reconciled F_2



- Least squares: Very poor DRPE
- Fair function: Better than least squares but inadequate
- Redescending: Avoids bounds on $1/A_1$ and $1/A_2$. Tuned for best result

Redescending Estimators superior for problems with largely corrupted data

Gross Errors Modeled with AIC

Common statistical framework for combinatorial and robust DRPE strategies?

- DR & Gross error detection \equiv model discrimination and PE
- Data models \Rightarrow partitions of residuals into random and gross errors
- Parameters = # of outliers + model parameters
- Obtain most likely model and its parameters
- Max. likelihood estimators (MLE): asymptotically efficient \Rightarrow likelihood function = sensitive criterion of deviation of parameters from true values
- Akaike Information Criterion (AIC): Distance between fitted and true models

$$AIC = -2 \cdot \log(\text{maximum likelihood}) + 2(\# \text{ independently adjusted parameters in model})$$

AIC for Data Reconciliation

Objective Function

$$AIC = -2\log(\text{maximum likelihood}) + 2(\# \text{ independently adjusted parameters within the model}) .$$

$$AIC = E(S) = 2 \sum_{i=1}^N -\log(\ell(\varepsilon(i, p), i, p)) + 2 \dim(p),$$

- Contains discrete and continuous variables
- dim(p) = # gross errors in problem
- Leads to MINLP problem
- Can be simplified to MILP for linear systems
- AIC can also be used as an off-line objective for calibration

83

MINLP formulation (Yamamura et al, 1988; Arora and B, 2001)

$$\begin{aligned} \min_{x_i, \mu_i, y_i, z_i} & \sum_{i=1}^n \left[(x_i^M - x_i) / \sigma_i - \mu_i / \sigma_i \right]^2 + 2 \cdot \sum_{i=1}^n y_i \\ \text{s.t.} & \boxed{Ax = 0} \\ & \mu_i \leq U_i y_i \\ & -\mu_i \leq U_i y_i \\ & \mu_i - z_i U_i - z_i L_i + L_i y_i \leq 0 \\ & -\mu_i + z_i U_i + z_i L_i + L_i y_i \leq L_i + U_i \\ & x_i \geq 0 \\ & z_i \leq y_i \\ & y_i, z_i \in \{0, 1\}. \end{aligned}$$

-Direct minimization of AIC function

-Model is for linear system, but straightforward extension for nonlinear systems

- μ_i = positive or negative biases,
 $U \geq |\mu_i| \geq L$

-Binary variables enforce this:

$y = 0, 1$: existence of gross error

$z = 0, 1$: positive or negative sign

84

MILP Simplification (Soderstrom et al., 2000)

$$\begin{aligned} \min_{x_i, \mu_i, y_i, z_i} & \sum_{i=1}^n \left| (x_i^M - x_i) / \sigma_i - \mu_i / \sigma_i \right| + \sum_{i=1}^n w_i \cdot y_i \\ \text{s.t.} & Ax = 0 \\ & \mu_i \leq U_i y_i \\ & -\mu_i \leq U_i y_i \\ & \mu_i - z_i U_i - z_i L_i + L_i y_i \leq 0 \\ & -\mu_i + z_i U_i + z_i L_i + L_i y_i \leq L_i + U_i \\ & z_i \leq y_i \\ & y_i, z_i \in \{0, 1\}, \end{aligned}$$

$$\begin{aligned} \min_{x_i, \mu_i, y_i, z_i} & \sum_{i=1}^n (r_i + q_i) / \sigma_i + \sum_{i=1}^n w_i \cdot y_i \\ \text{s.t.} & Ax = 0 \\ & x_i^M - (x_i + \mu_i) = r_i - q_i \\ & \mu_i \leq U_i y_i \\ & -\mu_i \leq U_i y_i \\ & \mu_i - z_i U_i - z_i L_i + L_i y_i \leq 0 \\ & -\mu_i + z_i U_i + z_i L_i + L_i y_i \leq L_i + U_i \\ & z_i \leq y_i \\ & r_i, q_i \geq 0 \\ & y_i, z_i \in \{0, 1\}. \end{aligned}$$

- Quadratic terms become absolute values (like Huber)
- Weighting values (w_i) not clear
- Limited to linear models
- Much faster solution times

85

Mixed Integer Formulation for Gross Errors

Minimize AIC by solving MINLP:

$$\begin{aligned} \min_{x_i, \mu_i, y_i} & \sum_{i=1}^n \left[\frac{(x_i^M - x_i)}{\sigma_i} - \frac{\mu_i}{\sigma_i} \right]^2 + 2 \cdot \sum_{i=1}^n y_i \\ \text{s.t.} & Ax = 0, \quad (-1) \\ & |\mu_i| \leq U_i y_i, \\ & |\mu_i| \geq L_i y_i, \\ & y_i \in \{0, 1\}, \\ & x_i \geq 0. \end{aligned}$$

MILP (Soderstrom et al., 2000):

$$\begin{aligned} \min_{x_i, \mu_i, y_i} & \sum_{i=1}^n \left| \frac{(x_i^M - x_i)}{\sigma_i} - \frac{\mu_i}{\sigma_i} \right| + \sum_{i=1}^n w_i \cdot y_i \\ \text{s.t.} & Ax = 0, \quad (-2) \\ & |\mu_i| \leq U_i y_i, \\ & |\mu_i| \geq L_i y_i, \\ & y_i \in \{0, 1\}, \\ & x_i \geq 0. \end{aligned}$$

w_i to obtain independent biases

- Both MINLP and MILP: Expensive for EVM problems
- In MILP: Choosing weights not entirely clear
- Neither suitable for large EVM problems

Cheaper + reliable DRPE in presence of outliers \Rightarrow use Robust Statistics

86

Tuning the Redescending Parameters

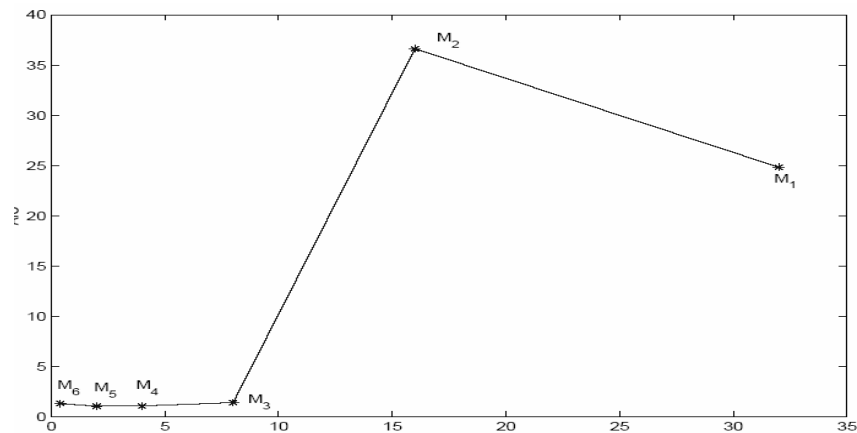
- Idea: *Minimize* AIC associated with the redescending estimator
- Tuning constants related to contamination in data
- Use size of redescending region to assign b and a to obtain c from $c = b + 2a$
- Simple Two-Step Tuning Procedure

- | | |
|--------------------------------------------|-----------|
| 1. Select range of redescending estimators | Step
1 |
| 2. Perform DRPE and obtain values of AIC | |

- | | |
|------------------------------------------|-----------|
| 3. Obtain bounds for location of minimum | Step
2 |
| 4. Golden section search for c | |
| 5. $b = c/2, a = (c - b)/2$ | |

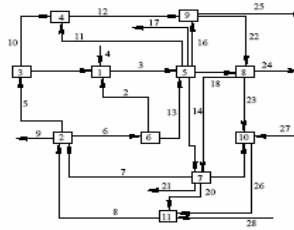
87

Tune Redescending Estimator Using AIC Tanks Example



88

Steam Metering Example



- 28 streams, 11 nodes \Rightarrow 11 linear equalities, 28 variables

- Moving horizon of data used $X^M = \begin{bmatrix} x_{11}^M & x_{12}^M & \dots & x_{1H}^M \\ x_{21}^M & x_{22}^M & \dots & x_{2H}^M \\ \vdots & \vdots & \ddots & \vdots \\ x_{n1}^M & x_{n2}^M & \dots & x_{nH}^M \end{bmatrix}$

- Sizes of horizon attempted: 5, 10, 20, 40
- 3, 5, 7 streams biased; bias propagation
- 100 moving horizons for fixed location of bias
- Smoothed redescending estimator for use in GAMS

89

Steam Metering Results

$$OP = \frac{\# \text{ biases correctly identified}}{\# \text{ biases simulated}}, \quad AVT1 = \frac{\# \text{ unbiased wrongly identified}}{\# \text{ simulation trials}}$$

Sample results for horizon size = 5, 3 biases:

Method	w_i	OP	AVT1	\mathcal{AIC}	CPU (s)/ horizon
Fair function	-	0.640	3.907	-	-
RE (a, b, c)					
0.5, 1, 2	-	0.673	5.820	3.491	1.114
1, 2, 4	-	0.626	2.065	14.112	1.075
2, 4, 8	-	0.571	0.305	40.299	1.057
4, 8, 16	-	0.455	0.069	96.399	1.019
MINLP	-	0.885	0.071	1.067	42.989
MILP	1	0.685	2.265	1.211	4.369
MILP	2H	0.704	0.081	1.117	0.592
MILP	100	0.112	0.000	17.148	0.386

- MIPs: used info. abt. type of gross error; EVM too expensive
- RE tailored to process \Rightarrow Tune RE parameters to minimize \mathcal{AIC}

90

Conclusions – Data Reconciliation

- Redescending estimators superior for gross error detection and provide good parameter estimates
- Redescending estimators
 - More robust than Huber estimator: Fair function
 - Can be tuned
- Simple two-step tuning strategy utilized: considerably improves DRPE
- MINLP and MILP: good but computationally intensive and useful only for linearly constrained problems